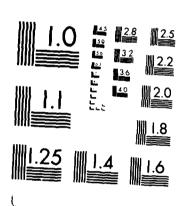
AD-R174 941 THE EFFECTS OF VARIANCE FUNCTION ESTINATION ON PREDICTION AND CALIBRATION (U) MORTH CAROLINA UNIV AT CHAPEL HILL DEPT OF STATISTICS R J CARROLL AUG 86 UNCLASSIFIED HIMEO-SER-1703 AFOSR-TR-86-2139 F/G 12/1 NL



CRGCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS 1963 A

AD-A174 941



"The Effects of Variance Function Estimation on Prediction and Calibration: An Example" 12. PERSONAL AUTHORIS) Carroll, Raymond J. 13b. TIME COVERED FROM 12		PORT DOCUME	NTATION PAGE	 E		
2. DISTRIBUTIONIZAVALABILITY OF REPORT 2. DECLASSIFICATION/DOWNGRADING SCHEDULE 2. PERFORMING ORGANIZATION AEPORT NUMBERIS) 3. NAME OF PERFORMING ORGANIZATION BID OFFICE SYMBOL IN APPROVED OF SCIENTIFIC RESEARCH 4. NONITORING ORGANIZATION REPORT NUMBERIS) 3. NAME OF PERFORMING ORGANIZATION BID OFFICE SYMBOL IN APPROVED OF SCIENTIFIC RESEARCH 4. NONITORING ORGANIZATION AT FORCE OF SCIENTIFIC RESEARCH 4. NONITORING ORGANIZATION AT FORCE OF SCIENTIFIC RESEARCH 4. NONITORING ORGANIZATION AT FORCE OF SCIENTIFIC RESEARCH 5. NONITORING ORGANIZATION NUMBER 5. NAME OF FUNDING/SPONGORING ORGANIZATION 5. NAME OF FUNDING/SPONGORING ORGANIZATION NUMBER 5. NAME OF FUNDING/SPONGORING ORGANIZATION NUMBER 6. NAME OF FUNDING/SPONGORING 6. NAME OF FUNDING ORGANIZATION NUMBER 6. NAME OF FUNDING ORGANI			1b. RESTRICTIVE MARKINGS			
AFOSR-TR- 86-2139 S. MONITORING ORGANIZATION REPORT NUMBERIS) AFOSR-TR- 86-2139 S. MONITORING ORGANIZATION REPORT NUMBERIS) AFOSR-TR- 86-2139 S. MONITORING ORGANIZATION AFOSR-TR- 86-2139 S. MONITORING ORGANIZATION AFOSR-TR- 86-2139 N. AME OF MONITORING ORGANIZATION AIF FORCE OFFICE SYMBOL IN AME OF MONITORING ORGANIZATION AIF FORCE OFFICE SYMBOL Philips Hall, Chapel Hill, NC 27514 N. ADDRESS (Cit), Sizer and ZIP Code; NAME OF FUNDING/SPONSORING ORGANIZATION NO.	28. SECURITY CLASSIFICATION AUTHORITY		Approved for public release;			
AFOSR-TR- 86-2139 6a NAME OF PERFORMING ORGANIZATION University of NO-Chapel Hill 6c. Address Interval 2IP Code: Univ. of North Carolina, Statistics Dept. Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Phillips Hall, Chapel Hill, NC 27514 6c. Address Interval 2IP Code: Balling Air Force Base Washington, DC 20332 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification: Phillips Hall, Chapel Hill, NC 27514 11. TITLE Interval Security Classification Phillips Hill, NC 27514 11. TITLE Interval Security Classification Phillips Hill, NC 27514 11. TITLE Interval Security Classification Phillips Hill, NC 27514 11. TITLE Interval Security Phillips Hill, NC 27514 11. TITLE Interval Security Phillips Hill, NC 27514 12. ABSTRACT Continue on Payers of Phillips Hill, NC 27514 13. ABSTRACT Continu	20. 00.000					
De la correction de la			[] [] [] [] [] [] [] [] [] []			
APPRESS (City, State and ZIP Code) No. ADDRESS (City, State and ZIP Code) Phillips Hall, Chapel Hill, NC 27514 La NAME OF FUNDING/SPONSORING OR AFOSR Re ADDRESS (City, State and ZIP Code) AFOSR Re ADDRESS (City, State and ZIP Code) Bolling Afr Force Base Washington, DC 20332 TITITLE (Include Security Classification) "The Effects of Variance Function Estimation on Prediction and Calibration: An Example" "Lee Personal Authorisis Carroll, Raymond J. 12 YERSONAL AUTHORIS Carroll, Raymond J. 13 YUFF OF REPORT technical FROM 8/85 to 8/86 August 1986 14 Date OF REPORT (Yr. Mo. Day) technical FROM SUB. GR. Weighted least squares; heteroscedasticity, regression 15 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 16 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 17 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 18 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 18 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 19 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 19 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 19 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 19 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 10 ASSTRACT Continue on reverse if receivery and identify by block number) Weighted least squares; heteroscedasticity, regression 10 ASSTRACT Continue on reverse if receivery an		66. OFFICE SYMBOL	78. NAME OF MONITORING ORGANIZATION			
Univ. of North Carolina, Statistics Dept. Phillips Hall, Chapel Hill, NC 27514 En NAME OF FUNDING/SPONSORING ORGANIZATION AFOSR AFOSR AFOSR AFOSR F-49620-85-C-0144. Re address (cir.) State and ZIP Code; Bolling Air Force Base Washington, DC 20332 11. TITLE Include Security Classification "The Effects of Variance Function Estimation on Prediction and Calibration: An Example" 12. PERSONAL AUTHORIST 13b. TIME COVERED ASS AUGUST 12 13c. TYPE OF REPORT FROM SIGN AUGUST 12 13c. TYPE OF REPORT FROM SUB. OR Weighted least squares; heteroscedasticity, regression 14c. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 15. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 16. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 17. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 18. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 18. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) Weighted least squares; heteroscedasticity, regression		(If applicable)	Air Force Office of Scientific Research			
Phillips Hall, Chapel Hill, NC 27514 Lande of Funding Sponsoring organization with a processory and contained in the processory and contained with a processor organization or prediction and Calibration: An Example with Embedding and Calibration and Calibration: An Example with Embedding and Calibration and Calibration and Calibration: An Example with Embedding and Calibration and Ca					_	_
AFOSR I SOURCE OF FUNDING NOS. PROGRAM ELEMENT NO. PROJECT NO. NO. ASSUMBLEMENT NO. PROJECT NO. ASSUMBLEMENT NO. PROJECT NO. ASSUMBLEMENT NO. PROJECT NO. ASSUMBLEMENT			same 45 8c			
Bolling Air Force Base Washington, DC 20332 TITLE Include Security Clausification PROGRAM PROJECT TASK NO.		(If applicable)				
Bolling Air Force Base Washington, DC 20332 11. TITLE Includes Security Classification "The Effects of Variance Function Estimation on Prediction and Calibration: An Example" 12. Personal Authoris; Carroll, Raymond J. 13a type of Report technical From 8/85 to 8/86 14 Date of Report (Yr. Mo. Dey) 15. PAGE COUNT technical From 8/85 to 8/86 14 Date of Report (Yr. Mo. Dey) 16. Supplementary Notation 17. Cosati codes weighted least squares; heteroscedasticity, regression 18. ABSTRACT Continue on reverse if necessary and identify by block number) - We consider, fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimater goal of a study may be a prediction or a calibration. Meishow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1 1.006 22b. TELEPHONE NUMBER 20. 22c. OFFICE SYMBIL. (Include Arms Code) 76-7	711 0311					
"The Effects of Variance Function Estimation on Prediction and Calibration: An Example" "The Effects of Variance Function Estimation on Prediction and Calibration: An Example" 12. PERSONAL AUTHORIS: Carroll, Raymond J. 13a TYPE OF REPORT technical FROM 8/85 TO 8/86 August 1986 12 16. SUPPLEMENTARY NOTATION 17. COSATI CODES 18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) - We consider fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Weishow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretica field that is not particularly well developed. 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1 1906 22b. TELEPHONE NUMBER 20 2c. OFFICE SYMBY. (Includes Area Coder) 7(-1) 22c. OFFICE SYMBY.	Bolling Air Force Base		PROGRAM	PROJECT		
12. PERSONAL AUTHORIS) Carroll, Raymond J. 13. TYPE OF REPORT technical FROM 8/85 TO 8/86 14. DATE OF REPORT (Yr., Mo., Dey) technical FROM 8/85 TO 8/86 12 15. SUPPLEMENTARY NOTATION 17. COSATI CODES 18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) FIELD GROUP SUB. GR. Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) - We consider, fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimater goal of a study may be a prediction or a calibration. Welshow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 21. ABSTRACT SECURITY CLASSIFICATION, DEC 1 1806 22. ABSTRACT SECURITY CLASSIFICATION, DEC 1 1806 23. ABSTRACT SECURITY CLASSIFICATION, DEC 1 1806 24. ABSTRACT SECURITY CLASSIFICATION, DEC 1 1806 25. ABSTRACT SECURITY CLASSIFICATION DEC 1 1806 26. ABSTRACT SECURITY SECURITY SECURITY CLASSIFICATION DEC 1 1806 26. ABSTRACT SECURITY SECURITY SECURITY SECURITY SECURITY SECURITY SECURITY SECURITY SECURITY SE	11. TITLE (Include Security Classification)		1			
Carroll, Raymond J. 13a. TIME COVERED		tion Estimation	lon Prediction	Land Calib	ration: An E	kample"
13b. TIME COVERED technical FROM 8/85 TO 8/86 August 1986 12 15. PAGE COUNT Lechnical FROM 8/85 TO 8/86 August 1986 12 16. SUPPLEMENTARY NOTATION 17. COSATI CODES FIELD GROUP SUB. GR. Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) We consider, fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Welshow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1 1906 22b. TELEPHONE NUMBER CO. 12c. OFFICE SYMBUL. (Include Area Code) 760						
18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) FIELD GROUP SUB.GR. Weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number) We consider fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Weishow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1 1500 NCLASSIFIED/UNLIMITED & SAME AS APT. DIICUSERS 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1 1500 NCLASSIFIED/UNLIMITED & SAME AS APT. DIICUSERS 220. TELEPHONE NUMBER 220. OFFICE SYMBER (Include Area Code; 767) 222. OFFICE SYMBER (Include Area Code; 767) 222. OFFICE SYMBER 223. OFFICE SYMBER 224. OFFICE SYMBER 224. OFFICE SYMBER 225. OFFICE SYMBER 226. OFFICE SYMBER						TNUC
weighted least squares; heteroscedasticity, regression 19. ABSTRACT Continue on reverse if necessary and identify by block number) We consider fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Weishow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 20. DISTRIBUTION/AVAILABILITY OF ABSTRACT NCLASSIFIED/UNLIMITED \$\mathbb{Z}\$ SAME AS APT. DIICUSERS \$\mathbb{L}\$ 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1.006 22b. TELEPHONE NUMBER \$\mathbb{Z}\$ 22c. OFFICE SYMBD. (Include Area Code: 767) 22c. OFFICE SYMBD.	والمستقد والمراجع والمراجن والمراجع والمراع والمراجع والمراع والمر					
weighted least squares; heteroscedasticity, regression 19. ABSTRACT Continue on reverse if necessary and identify by block number) We consider fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Weishow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 20. DISTRIBUTION/AVAILABILITY OF ABSTRACT NCLASSIFIED/UNLIMITED \$\mathbb{Z}\$ SAME AS APT. DIICUSERS \$\mathbb{Z}\$ TELEPHONE NUMBER \$\mathbb{Z}\$ DICT. DICT						
weighted least squares; heteroscedasticity, regression 19. ABSTRACT (Continue on reverse if necessary and identify by block number)	}	18. SUBJECT TERMS (C	Continue on reverse if necessary and identify by block number)			
We consider fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Welshow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1 1906 22. DISTRIBUTION/AVAILABILITY OF ABSTRACT NCLASSIFIED/UNLIMITED & SAME AS RPT. DICCUSERS DICCUSE	FIELD GROUP SUB. GR.	weighted leas	st squares; heteroscedasticity, regression			
In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Welshow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed. 21. ABSTRACT SECURITY CLASSIFICATION. DEC 1 1906 22. NAME OF RESPONSIBLE INDIVIDUAL (Include Area Code) 767 22. OFFICE SYMBOL (Include Area Code) 767	19. ABSTRACT (Continue on reverse if necessary and identify by block number)					
Q20. NAME OF RESPONSIBLE INDIVIDUAL Control	In most fields, it is fairly common folklore thathow one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. Weishow by an example that how one handles the variance function can have large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical					
226. NAME OF RESPONSIBLE INDIVIDUAL 226. TELEPHONE NUMBER 226. OFFICE SYMBOT. (Include Area Code) 767	20. DISTRIBUTION/AVAILABILITY OF ABSTRA	ст	21. ABSTRACT SEC	URITY CLASSIFI	CATHON DEC	1 1986 _{1 1}
(Include Area Code) 767	PINCLASSIFIED/UNLIMITED TO SAME AS RPT. TO DTIC USERS					
	litas C	1.10.	(Include Area Co	oder 767		BOT

DD FORM 1473, 83 APR

AFOSR-TR. 86-2139

THE EFFECTS OF VARIANCE FUNCTION ESTIMATION ON PREDICTION AND CALIBRATION : AN EXAMPLE

Raymond J. Carroll

University of North Carolina at Chapel Hill

Research supported by the Air Force Office of Scientific Research Contract

15007 F-49620-85-C-0144.

KEY WORDS AND PHRASES : Weighted least squares. Heteroscedasticity, Regression.

Accession For					
NTIS	GRA&I	X			
DTIC '	TAB				
Unannounced 🔲					
Justification					
Ву					
Distribution/					
Availability Codes					
Arail and/or					
Dist	Special				
1	1				
11					
HT					



ABSTRACT

We consider fitting a straight line to data when the variances are not constant. In most fields, it is fairly common folklore that how one estimates the variances does not matter too much when estimating the regression function. While this may be true, most problems do not stop with estimating the slope and intercept. Indeed, the ultimate goal of a study may be a prediction or a calibration. We show by an example that how one handles the variance function can has large effects. The point is almost trivial, but so often ignored that it is worth documenting. Additionally, this points out that one ought to spend time trying to understand the structure of the variability, a theoretical field that is not particularly well developed.

1 : INTRODUCTION

Consider a heteroscedastic regression model, in which we observe $\,N\,$ pairs (y,x) following the model

$$(1.1) E(y|x) = f(x,\beta) ;$$

(1.2) Standard Deviation(
$$y|x$$
) = $\sigma g(x, \beta, \theta)$.

While our remarks hold generally, in what follows it suffices to consider the special case of linear regression for the mean and the power of the mean model for the standard deviation, i.e.,

(1.3)
$$f(x,\beta) = \beta_0 + \beta_1 x : g(x,\beta,\theta) = f(x,\beta)^{\theta}.$$

When $\theta=0$, we have the homoscedastic regression model, and unweighted least squares will ordinarily be used to estimate β . For other values of θ , generalized least squares can be used to estimate β , see Carroll & Ruppert (1987) for a discussion and a review of the literature. Generalized least squares is weighted least squares with estimated weights. The version of generalized least squares used here for each θ is fully iterated reweighted least squares, sometimes called quasi-likelihood, see McCullagh & Nelder (1983). In practice, θ is unknown and must be estimated. The theory of such estimation is given by Davidian & Carroll (1986).

The common folklore theorem of generalized least squares states that as long as one's estimate $\hat{\theta}$ of θ is root-N consistent, the resulting generalized least squares estimate has the same asymptotic distribution as if θ were known.

See Judge, et al (1985) and Carroll & Ruppert (1982, 1987) for references and proofs. Indeed, any generalized least squares estimate has the same limit distribution as weighted least squares based on the correct weights, i.e., the inverse of the square of (1.2).

The folklore theorem has an analogue in practice. In the linear regression model with a reasonably sized data set, since unweighted least squares is consistent its fitted values rarely differ much from the fitted values from a generalized least squares fit. Consequently, the usual practice is to treat the estimation of the variance function $g(x,\beta,\theta)$ fairly cavalierly. if at all. To quote Schwartz (1979), "there is one point of agreement among statistics texts and that is the minimal effect of weighting factors on fitted regression curves. Unless the variance nonuniformity is quite severe, the curve fitted to calibration data is likely to be nearly the same, whether or not the variance nonuniforming is included in the weighting factors". narrow focus on estimating the mean is misplaced, as Schwartz later notes, see also Garden, et al (1980). Sometimes the variance function is itself of importance. Box & Meyer (1985) state that "one distinctive feature of Japanese quality control improvement techniques is the use of statistical experimental design to study the effect of a number of factors on variance as well as the Other times the variance function essentially determines the quantity of interest. This occurs, for example, in the estimation of the sensitivity of a chemical or biochemical assay, see Carroll, Davidian & Smith (1986). However, there are even more basic problems where the variance function is of considerable importance, namely prediction and calibration.

It is perhaps trite to state that how well one estimates the variance function has a large effect on how well one can do prediction and calibration. It is, however, a point that is rarely taken into account in practice, as any

review of the (rudimentary) techniques in the assay literature will quickly show. There are two ways to see this point. The first is through an asymptotic theory outlined in section 3, where we show that the difference in the length of a prediction interval between θ known and unknown is asymptotically distributed with variance a monotone function of how well one estimates θ . The second and probably more useful way to see the effect of variance function estimation is through an example. The large costs involved in not weighting at all will be evident in this example, and will serve as an object lesson.

2. : CALIBRATION AND PREDICTION

Calibration experiments start with a training or calibration sample $(y_1,x_1),\ldots,(y_N,x_N)$ and then fit models to the mean and variance structures. The real interest lies in an independent pair (y_0,x_0) . Sometimes x_0 is known and we wish to obtain confidence intervals for y_0 ; this is prediction. Other times. y_0 is easily measured but x_0 is unknown and inference is to be made about it, see Rosenblatt & Spiegelman (1982).

For example, in an assay x might represent the concentration of a substance and y might represent a counted value or intensity level which varies with concentration. One will have a new value y_{\square} of the count or intensity and wish to draw inference about the true concentration x_{\square} . The calibration sample is drawn so that we have a good understanding of how the reponse varies as a function of concentration. The regression equation relating the response to

concentration is then inverted to predict the concentration from the observed response.

For the remainder of this section we will assume that the responses are normally distributed, although this can be relaxed. Given a value \mathbf{x}_{\square} , the standard point estimate of the response \mathbf{y}_{\square} is $\mathbf{f}(\mathbf{x}_{\square}, \boldsymbol{\beta})$. Let $\hat{\boldsymbol{\beta}}_{\mathbf{G}}$ be a generalized least squares estimate, and define

$$S_{G} = S_{G}(\theta) = N^{-1} \sum_{i=1}^{N} f_{\beta}(x_{i}.\hat{\beta}_{G}) f_{\beta}(x_{i}.\hat{\beta}_{G})^{T} / f(x_{i}.\hat{\beta}_{G})^{2\theta} .$$

$$\hat{\sigma}^{2}(\theta) = N^{-1} \sum_{i=1}^{N} (y_{i} - f(x_{i}.\hat{\beta}_{G}))^{2} / f(x_{i}.\hat{\beta}_{G})^{2\theta} .$$

where f_{eta} is the derivative of f with respect to eta. For large calibration data sets, the variance in the error made by prediction is

(2.1)
$$\text{Variance}\{y_{\square} - f(x_{\square}, \hat{\beta}_{G})\} \cong \sigma^{2} q_{N}^{2}(x_{\square}, \beta, \theta), \text{ where }$$

$$q_{N}^{2}(x_{\square},\beta,\theta) = g^{2}(x_{\square},\beta,\theta) + N^{-1}f_{\beta}(x_{\square},\beta)^{T} S_{G}^{-1} f_{\beta}(x_{\square},\beta).$$

Note that if the size N of the calibration data set is large, then the error in prediction is determined predominately by the variance function

$$\sigma^2 g^2(f(x_{\square},\beta),x_{\square},\theta),$$

and not by the calibration data set itself. An approximate $(1-\alpha)100\%$ confidence interval for the response y_Π is given by

(2.2)
$$I(x_{\square}) = \{\text{all values y in the interval} \\ f(x_{\square}, \hat{\beta}_{G}) \pm t_{1-\alpha/2}^{N-p} \hat{\sigma}_{G} q_{N}(x_{\square}, \hat{\beta}_{G}, \theta) \} ,$$

where $t_{1-\alpha/2}^{N-p}$ is the $(1-\alpha/2)$ percentage point of a t-distribution with N-p degrees of freedom. For large sample sizes, this interval becomes

(2.3)
$$\begin{split} I(x_{\square}) &\cong \{\text{all y in the interval} \\ f(x_{\square}, \hat{\beta}_G) &\pm t_{1-\alpha/2}^{N-p} \hat{\sigma}_G \ g(f(x_{\square}, \hat{\beta}_G), x_{\text{new}}, \theta)\} \end{split} .$$

The prediction interval (2.2) is only an approximate $(1-\alpha)100\%$ confidence interval because the function $q_{_{\rm N}}$ is not known but rather must be estimated.

The effect of ignoring the heterogeneity can be seen through examination of (2.3). If $\hat{\sigma}_L^{-2}$ is the unweighted mean squared error, then for large samples we have the approximation

$$\hat{\sigma}_{L}^{2} \cong \sigma^{2} g_{\text{mean}}^{2} = \sigma^{2} N^{-1} \Sigma_{i=1}^{N} g^{2}(x_{i},\beta,\theta)$$

Thus the unweighted prediction interval for large sample sizes is approximately

(2.4)
$$I_{L}(x_{\square}) \cong \{ \text{ all y in the interval} \\ f(x_{\square}, \hat{\beta}_{L}) \pm t_{1-\alpha/2}^{N-p} \sigma g_{\text{mean}} \} .$$

Comparing (2.3) and (2.4) we see that where the variability is small, the unweighted prediction interval will be too long and hence pessimistic, and conversely where the variance is large.

Now suppose that we are given the value of the response y_{\square} and wish to estimate and make inference about the unknown x_{\square} . The estimate of x_{\square} is that

value which satisfies $f(x_{\square}, \beta_{G}) = y_{\square}$. The most common interval estimate of x_{\square} is the set of all values x for which y_{\square} falls in the prediction interval I(x), i.e.

Calibration interval for
$$x_{\square} = \{ \text{ all } x \text{ such that } y_{\square} \in I(x) \}$$

where $I(x)$ is given by (2.3)

The effect of not weighting is too long and pessimistic confidence intervals for \mathbf{x}_{\square} where the variance is small and the opposite where the variance is large. As far as we know, little work has been done to determine whether one can shorten the calibration confidence interval by making more direct use of the variance function.

3. : ASYMPTOTICS

Assume throughout that the data are symmetrically distributed about their mean. Let $\hat{\beta}_{\rm G}$ be any generalized least squares estimate of β based on an estimate of θ , call it $\hat{\theta}$ say. Davidian & Carroll (1986) introduce a class of estimators which depend on the data only through $\hat{\beta}_{\rm G}$, the design $\{{\bf x}_i\}$, and either sample variances from replicates at each design point or on transformations of the squared residuals

$$\{y_i - f(x_i, \hat{\beta}_G)\}^2$$

This class of estimators includes most methods in the literature, see Judge, et

al (1985). Davidian & Carroll (1986) show that all members of their class of estimators have the asymptotic expansion

$$N^{1/2}(\hat{\theta} - \theta) = W_N + a^T N^{1/2}(\hat{\beta}_G - \beta) + o_p(1) .$$

In (3.1), a is a fixed vector and \mathbf{W}_N is asymptotically normally distributed. Because the observations have symmetric distribution, \mathbf{W}_N is asmptotically uncorrelated with $\hat{\boldsymbol{\beta}}_G$.

Let $\hat{\beta}_{G}(\theta)$ and $\hat{\beta}_{G}(\hat{\theta})$ be generalized least squares estimates of β with θ known and unknown respectively, and let $\hat{\sigma}(\theta)$ and $\hat{\sigma}(\hat{\theta})$ be the corresponding estimates of σ . The length of the prediction intervals with θ known and unknown are proportional to $L(\theta)$ and $L(\hat{\theta})$ respectively, where

$$L(\theta) = \hat{\sigma}(\theta) q_N(x_0, \hat{\beta}_G(\theta), \theta)$$
.

The random variable

$$\Delta L = N^{1/2} \{ L(\hat{\theta}) - L(\theta) \} / \sigma$$

describes how well one approximates the length one would use if θ were known. Intuitively, we would like ΔL to have smallest possible variability.

THEOREM: Suppose that W_N in (3.1) is asymptotically normally distributed with mean zero and covariance $\mathbf{C} = \mathbf{C}(\hat{\boldsymbol{\theta}})$ depending on the method of estimating $\boldsymbol{\theta}$. Then, under regularity conditions, $\Delta \mathbf{L}$ is asymptotically normally distributed with variance an increasing function of $\mathbf{C}(\hat{\boldsymbol{\theta}})$.

NOTE: The Theorem remains valid if ΔL is the normalized difference in length between the interval with θ unknown and the interval with completely specified variance function.

<u>PROOF</u> (Sketch): It is easily seen that in the definition of L we may replace q_N by g. Further, $\Delta L = A_1 + A_2 + A_3$, where

$$\begin{split} \mathbf{A}_1 &= \mathbf{N}^{1/2} \mathbf{g}(\mathbf{x}_{\square}, \hat{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}}), \hat{\boldsymbol{\theta}}) \ \{ \hat{\boldsymbol{\sigma}}(\hat{\boldsymbol{\theta}}) - \hat{\boldsymbol{\sigma}}(\boldsymbol{\theta}) \} / \boldsymbol{\sigma} \\ \mathbf{A}_2 &= \mathbf{N}^{1/2} \ \{ \hat{\boldsymbol{\sigma}}(\boldsymbol{\theta}) / \boldsymbol{\sigma} \} \ \{ \mathbf{g}(\mathbf{x}_{\square}, \hat{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}}), \hat{\boldsymbol{\theta}}) - \mathbf{g}(\mathbf{x}_{\square}, \hat{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}}), \boldsymbol{\theta}) \} \\ \mathbf{A}_3 &= \mathbf{N}^{1/2} \{ \hat{\boldsymbol{\sigma}}(\boldsymbol{\theta}) / \boldsymbol{\sigma} \} \ \{ \mathbf{g}(\mathbf{x}_{\square}, \hat{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}}), \boldsymbol{\theta}) - \mathbf{g}(\mathbf{x}_{\square}, \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}) \} \end{split}$$

Now. $A_3 \xrightarrow{p} 0$ since, from Carroll & Ruppert (1987), we have that

$$N^{1/2}\{\hat{\beta}(\hat{\theta}) - \hat{\beta}(\theta)\} / \sigma \xrightarrow{p} 0,$$

By a Taylor series.

$$A_2 = N^{1/2} g_{\chi}(x_{\square}, \beta, \theta) (\hat{\theta} - \theta) + o_{p}(1) .$$

Lemma A.3 of Carroll, Davidian & Smith shows that for some constant b(x $_{\square}$),

$$A_1 = b(x_0) N^{1/2} (\hat{\theta} - \theta) + o_p(1)$$
.

This shows that for some constant $c(x_{\underline{u}})$.

$$\Delta L = c(x_0) N^{1/2} (\hat{\theta} - \theta) + o_p(1)$$

The proof is completed by applying (3.1).

4. : AN EXAMPLE

In Chapter 2, section 8, Carroll & Ruppert (1987) present the results of an assay for the concentration of an enzyme (esterase). There were 113 observations, of which 5 were deleted. The observed concentration of esterase was recorded and then a binding experiment was undertaken, so that the response is the count of the number of bindings. These data were given to us by another statistician and we are unable to give further detail into the background of the experiment. We do not know wnether the recorded concentration of esterase has been accurately measured, although we will assume it has been and that there is little if any measurement error in this predictor. The lack of replicates in the reponse is rather unusual in our experience. Since the response is a count, one might expect Poisson variation, i.e., the power of the mean model holds with $\theta = 0.50$. In our experience with assays, such a model almost always underestimates θ , with values between 0.60 and 0.90 being much more common: see Finney (1976) and Raab (1981a).

The eventual goal of the study is to take observed counts and infer the concentration of esterase, especially for smaller values of the latter. As is typical in these experiments, a calibration or training data set is taken for which the predictor variable esterase is known as is the counted response. Carroll & Ruppert (1982) plot the data, which appears reasonably although not perfectly linear. Actually, the logarithm of the response plotted against the

logarithm of the predictor may appear more linear to some, and less heteroscedastic. As in evident from that plot, the data exhibit rather severe heterogeneity of variance. The Spearman correlation between absolute studentized residuals and predicted values from an unweighted least squares fit is $\rho=0.39$ with formal computed significance level ≤ 0.0001 . Analysis as in Carroll & Ruppert (1982) indicate that the constant coefficient of variation model $\theta=1.0$ is reasonable, although a value $\theta=0.9$ might be even better. For $\theta=1.0$, the Spearman correlation between absolute studentized residuals and predicted values is $\rho=-0.10$, with significance level 0.29. In Figure 1, we plot kernel regression estimates of the Anscombe studentized residuals. i.e.. the absolute studentized residuals to the power 2/3, see McCullagh & Nelder (1983). Note that the plots indicate that $\theta=1.0$ does a far better job of accounting for the heteroscedasticity.

the states course wants and analysis

In these data, the effect of not weighting should be to have prediction and calibration confidence intervals which are much too large for small amounts of esterase and conversely for large amounts. In Figure 2 we plot the 95% prediction intervals for the count response for unweighted versus weighted regression: the effect is clear. A similar plot for the calibration intervals shows the same effect: the unweighted analysis is much too conservative for small amounts of esterase, and much too liberal for larger amounts. As Oppenheimer, et al (1983) state, "Rather dramatic differences have been observed depending on whether a valid weighted or invalid unweighted analysis is used".

This example shows that the actual prediction intervals are sensitive to misspecification of the variance function. It should be clear by inference and the previous section that one should make efforts to estimate the structural variance parameter θ as well as possible.

REFERENCES

- Box, G. E. P. & Myer, R. D. (1985). Dispersion effects from fractional designs. <u>Technometrics</u>, 28, 19-28.
- Carroll, R. J., Davidian, M. & Smith, W. (1986). Variance function estimation and the minimum detectable concentration in assays. Preprint.
- Carroll, R. J., and Ruppert, D. (1982). Robust estimation in teroscedastic linear models, <u>Annals of Statistics</u> 10, 429-441.
- Carroll, R. J., and Ruppert, D. (1987). <u>Transformations and Weighting in Regression</u>. Chapman & Hall, London.
- Davidian, M. and Carroll, R.J. (1986) Variance function estimation in regression. Preprint.
- Finney, D. J. (1976). Radioligand assay. Biometrics 32, 721-740.
- Garden, J. S., Mitchell, D. G. & Mills, W. N. (1980). Nonconstant variance regression techniques for calibration curve based analysis. Analytical Chemistry 52, 2310-2315.
- McCullagh, P. & Nelder, J. A. (1983). <u>Generalized Linear Models</u>. Chapman & Hall, New York.

session artificial presents represents additional contracts

- Oppenheimer, L., Capizzi, T.P., Weppelman, R.M. and Mehto, H. (1983)
 Determining the lowest limit of reliable assay measurement.

 <u>Analytical Chemistry</u> 55, 638-643.
- Raab, G. M. (1981). Estimation of a variance function, with application to radioimmunoassay. <u>Applied Statistics</u> 30, 32-40.
- Rosenblatt, J. R. & Spiegelman, C. H. (1981). Discussion of the paper by Hunter & Lamboy. <u>Technometrics</u> 23, 329-333.
- Schwartz, L. M. (1979). Calibration curves with nonuniform variance.

 <u>Analytical Chemistry</u> 51, 723-729.

FIGURE 1

The esterase assay data. This is a plot of the kernel regression fits to the Anscombe absolute residuals against the logarithms of the predicted values. The unweighted least squares fit is the solid line, while the generalized least squares fit for the constant coefficient of variation model is the dashed line. Endpoint effects have been ajusted for by selective deletion.

Figure 2

The esterase assay data. These are the 95% prediction intervals for a new response. The dashed line is unweighted least squares, while the solid line is the constant coefficient of variation fit. The lower part of the least squares interval has been truncated at zero where necessary.